

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Application of AOKI, *et al.*

Application No.	Unassigned	Examiner:	Unassigned
Filed:	Herewith	Group Art Unit:	Unassigned
For:	MORPHOLOGICAL ANALYZER, NATURAL LANGUAGE PROCESSOR, MORPHOLOGICAL ANALYSIS METHOD AND PROGRAM		

CLAIM OF FOREIGN PRIORITY

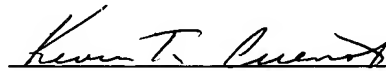
Box Patent Application
Commissioner for Patents
P.O. Box 1450
Alexandria, VA 22313-1450

Sir:

Priority under the International Convention for the Protection of Industrial Property and under 35 U.S.C. §119 is hereby claimed for the above-identified patent application, based upon Japanese Patent Application No. 2003-033220 filed February 12, 2003, and a certified copy of this application is submitted herewith which perfects the Claim of Foreign Priority.

Respectfully submitted,

Date: 2/12/04



Gregory A. Nelson, Registration No. 30,577
Kevin T. Cuenot, Registration No. 46,283
Brian K. Buchheit, Registration No. 52,667
AKERMAN SENTERFITT
Post Office Box 3188
West Palm Beach, FL 33402-3188
Telephone: (561) 653-5000

Express Mailing Label No. EV 347797285 US

日本国特許庁
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出願年月日

Date of Application:

2003年 2月12日

出願番号

Application Number:

特願2003-033220

[ST.10/C]:

[JP2003-033220]

出願人

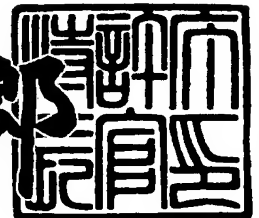
Applicant(s):

インターナショナル・ビジネス・マシーンズ・コーポレーション

2003年 5月 9日

特許庁長官
Commissioner,
Japan Patent Office

太田信一郎



出証番号 出証特2003-3034932

【書類名】 特許願

【整理番号】 JP9020244

【提出日】 平成15年 2月12日

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 17/27

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ピー・エム株式会社 大和事業所内

【氏名】 青木 和夫

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ピー・エム株式会社 大和事業所内

【氏名】 井ノ川 浩

【発明者】

【住所又は居所】 神奈川県大和市下鶴間 1 6 2 3 番地 1 4 日本アイ・ピー・エム株式会社 大和事業所内

【氏名】 中山 章弘

【特許出願人】

【識別番号】 390009531

【氏名又は名称】 インターナショナル・ビジネス・マシーンズ・コーポレーション

【代理人】

【識別番号】 100086243

【弁理士】

【氏名又は名称】 坂口 博

【代理人】

【識別番号】 100091568

【弁理士】

【氏名又は名称】 市位 嘉宏

【代理人】

【識別番号】 100108501

【弁理士】

【氏名又は名称】 上野 剛史

【復代理人】

【識別番号】 100104880

【弁理士】

【氏名又は名称】 古部 次郎

【選任した復代理人】

【識別番号】 100118201

【弁理士】

【氏名又は名称】 千田 武

【手数料の表示】

【予納台帳番号】 081504

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9706050

【包括委任状番号】 9704733

【包括委任状番号】 0207860

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 形態素解析装置、自然言語処理装置、形態素解析方法及びプログラム

【特許請求の範囲】

【請求項 1】 処理対象の自然言語文に対して形態素解析を行う形態素解析装置において、

見出し語および当該見出し語の属性情報を格納した辞書部と、

前記辞書部を参照して、処理対象の前記自然言語文から当該自然言語文を構成することが可能なトークンを抽出し、トークンリストに登録するトークンリスト作成部と、

前記トークンリスト作成部にて作成されたトークンリストに基づいて前記自然言語文を構成するのに最適なトークン列を選択するトークン列選択部とを備え、

前記トークンリスト作成部は、形態素解析に対して与えられた条件と、抽出した前記トークンに対応する前記見出し語の属性情報とに基づいて、当該トークンの前記トークンリストへの登録制御を行うことを特徴とする形態素解析装置。

【請求項 2】 前記トークンリスト作成部は、前記トークンに対応する前記見出し語の属性情報に基づき、形態素解析に対して与えられた前記条件に合致する属性を持つトークンのみを前記トークンリストに登録することを特徴とする請求項 1 に記載の形態素解析装置。

【請求項 3】 前記辞書部は、前記見出し語の属性情報として当該見出し語が分割可能か否かを示す情報を格納し、

前記トークンリスト作成部は、複合語を分割して形態素解析を行うという条件が与えられた場合に、前記見出し語の属性情報を参照し、抽出した前記トークンから分割可能な見出し語に対応するトークンを除いて前記トークンリストに登録することを特徴とする請求項 1 に記載の形態素解析装置。

【請求項 4】 前記辞書部に格納された前記見出し語の属性情報は、当該属性情報の数に応じたビット数のフラグデータで記録され、

前記トークンリスト作成部は、前記トークンに対応する前記見出し語の前記フラグデータの値に基づいて、当該トークンを前記トークンリストに登録するか否

かを決定することを特徴とする請求項 1 に記載の形態素解析装置。

【請求項 5】 処理対象の自然言語文に対して形態素解析を行う形態素解析装置において、

処理対象の前記自然言語文を当該自然言語文の構成要素であるトークンに分解し、より小さいトークンに分割可能なトークンを除いてトークンリストに登録するトークンリスト作成手段と、

前記トークンリスト作成手段にて作成されたトークンリストに基づいて前記自然言語文を構成するのに最適なトークン列を選択するトークン列選択手段とを備えることを特徴とする形態素解析装置。

【請求項 6】 前記トークンリスト作成手段は、形態素解析に与えられた所定の条件に応じて、前記トークンリストに登録するトークンからより小さいトークンに分割可能なトークンを除くか否かを選択的に制御可能であることを特徴とする請求項 5 に記載の形態素解析装置。

【請求項 7】 処理対象である自然言語文を形態素解析する形態素解析手段と、

前記形態素解析手段により形態素解析された前記自然言語文に対して所定の処理を行うアプリケーション実行手段とを備え、

前記形態素解析手段は、

見出し語および当該見出し語の属性情報を格納した辞書部と、

前記辞書部を参照して、処理対象の前記自然言語文から当該自然言語文を構成することが可能なトークンを抽出し、当該トークンに対応する前記見出し語の属性情報に基づいて、前記アプリケーション実行手段にて要求される条件に合致する属性を持つトークンのみをトークンリストに登録するトークンリスト作成部と

前記トークンリスト作成部にて作成されたトークンリストに基づいて前記自然言語文を構成するのに最適なトークン列を選択するトークン列選択部とを備えることを特徴とする自然言語処理装置。

【請求項 8】 前記辞書部は、前記見出し語の属性情報として当該見出し語が分割可能か否かを示す情報を格納し、

前記トークンリスト作成部は、分割可能な単語を分割して形態素解析を行うことが前記アプリケーション実行手段により要求される場合に、前記見出し語の属性情報を参照し、分割不可能な見出し語に対応するトークンを前記トークンリストに登録することを特徴とする請求項 7 に記載の自然言語処理装置。

【請求項 9】 前記辞書部に格納された前記見出し語の属性情報は、当該属性情報の数に応じたビット数のフラグデータで記録され、

前記トークンリスト作成部は、前記トークンに対応する前記見出し語の前記フラグデータの値に基づいて、当該トークンを前記トークンリストに登録するか否かを決定することを特徴とする請求項 7 に記載の自然言語処理装置。

【請求項 10】 コンピュータを用いて自然言語文に対し形態素解析を行う形態素解析方法であって、

処理対象の自然言語文を入力し、メモリに格納された辞書を参照して、当該自然言語文を構成することが可能なトークン及び当該トークンの属性情報を取得し、メモリの作業領域に格納するステップと、

形態素解析に与えられた所定の条件および前記トークンの属性情報に基づき、前記メモリに格納されたトークンの中から所定のトークンを選択してメモリの所定領域に構築されたトークンリストに登録するステップと、

前記トークンリストに基づいて処理対象の前記自然言語文を構成することが可能なトークン列を生成し、メモリの作業領域に格納するステップと、

前記メモリに格納された前記トークン列の中で処理対象の前記自然言語文を構成するのに最適なトークン列を選択し出力するステップとを含むことを特徴とする形態素解析方法。

【請求項 11】 前記トークンを前記トークンリストに登録するステップでは、前記形態素解析に与えられた所定の条件に応じて、当該所定の条件に合致する属性を持つトークンのみを前記トークンリストに登録することを特徴とする請求項 10 に記載の形態素解析方法。

【請求項 12】 コンピュータを用いて自然言語文に対し形態素解析を行う形態素解析方法であって、

処理対象の自然言語文を入力し、当該自然言語文の構成要素であるトークンに

分解し、得られたトークン群をメモリの作業領域に格納するステップと、

前記トークン群を、より小さいトークンに分割可能なトークンを除いてメモリの所定領域に構築されたトークンリストに登録するステップと、

前記トークンリストに基づいて処理対象の前記自然言語文を構成することが可能なトークン列を生成し、メモリの作業領域に格納するステップと、

前記メモリに格納された前記トークン列の中で処理対象の前記自然言語文を構成するのに最適なトークン列を選択し出力するステップと

を含むことを特徴とする形態素解析方法。

【請求項 1 3】 コンピュータを制御して、自然言語文の形態素解析を行うプログラムにおいて、

所定の記憶装置に格納され見出し語および当該見出し語の属性情報を記録した辞書を参照して、処理対象の前記自然言語文から当該自然言語文を構成することが可能なトークンを抽出し、形態素解析に与えられた所定の条件および前記トークンの属性情報に基づき、抽出されたトークンから所定のトークンを選択してメモリの所定領域に構築されたトークンリストに登録する手段と、

前記トークンリスト作成部にて作成されたトークンリストに基づいて前記自然言語文を構成するのに最適なトークン列を選択する手段として前記コンピュータを機能させることを特徴とするプログラム。

【請求項 1 4】 前記トークンを前記トークンリストに登録する手段では、前記辞書に記録された前記トークンの属性情報の数に応じたビット数のフラグデータの値に基づいて、当該トークンを前記トークンリストに登録するか否かを決定することを特徴とする請求項 1 3 に記載のプログラム。

【請求項 1 5】 コンピュータを制御して、自然言語文の形態素解析を行うプログラムにおいて、

処理対象の自然言語文を入力し、当該自然言語文の構成要素であるトークンに分解し、得られたトークン群をメモリの作業領域に格納する第 1 の処理と、

前記トークン群を、より小さいトークンに分割可能なトークンを除いてメモリの所定領域に構築されたトークンリストに登録する第 2 の処理と、

前記トークンリストに基づいて処理対象の前記自然言語文を構成することが可

能なトークン列を生成し、メモリの作業領域に格納する第3の処理と、

前記メモリに格納された前記トークン列の中で処理対象の前記自然言語文を構成するのに最適なトークン列を選択し出力する第4の処理とを前記コンピュータに実行させることを特徴とするプログラム。

【請求項16】 形態素解析に与えられた所定の条件を判断する処理を前記コンピュータにさらに実行させ、

前記所定の条件に応じて、前記第2の処理に代えて、全ての前記トークンを前記トークンリストに登録する処理を前記コンピュータに実行させることを特徴とする請求項15に記載のプログラム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、コンピュータを用いた自然言語解析に関し、特に形態素解析等の文章を単語に分解する技術に関する。

【0002】

【従来の技術】

コンピュータを用いた自然言語解析では、まず文章を単語に分解することが行われる。日本語などのように単語を区切らない表記方法を採用する言語では、形態素解析を行って文章を構成する単語が抽出される。

このように文章を単語に分解する処理では、2つ以上の単語が結びついて1つの単語を形成している複合語を適切に分割することが重要であり、従来から種々の技術が存在する（例えば、特許文献1参照）。

【0003】

図11は、コンピュータにて実現される従来の形態素解析手段の機能ブロックを示す図、図12は、従来の形態素解析のアルゴリズムを概略的に説明するフローチャートである。

図11、12に示すように、形態素解析では、まずトークンリスト作成部111が、処理対象の文章から様々なサイズの文字列を切り出し、全ての可能性のあるトークンを得る（ステップ1201）。そして、マスター辞書112を検索し

、各トークンとその属性（品詞など）を登録したトークンリストを作成する（ステップ1202）。ここで、トークンとは、文章や単語を構成する最小の要素であり、例えば「形態素」という語では、「形」、「形態」、「形態素」、「態」、「素」がそれぞれトークンとなる。

【0004】

次に、トークン列選択部113が、文法辞書114を参照し、ステップ1201で検出された全ての可能性のあるトークンの組み合わせの中から最も適切なトークン列を選ぶ（ステップ1203）。

この後、複合語分割処理部115が、ステップ1203で選択されたトークン列に対して複合語辞書116とのマッチングを行い、分割可能なトークンを更に細かいトークンに分割する（ステップ1204）。

【0005】

【特許文献1】

特開2002-251402号公報

【0006】

【発明が解決しようとする課題】

上述したように従来形態素解析では、トークン列を選択した後に複合語の分割処理を行っているため、複合語の部分に対するマッチング処理等の分だけ別途に時間を要し、この時間は文章中に含まれる複合語が多いほど長くなっていた。

また、上記従来形態素解析では、適切なトークン列を選択した後に複合語の分割処理を行っているため、分割された単語（トークン）によるトークン列が最適なものかどうか保証されないという欠点があった。

さらに、複合語の分割処理において参照される複合語辞書は、複合語及び複合語を構成する各単語に関して品詞情報や区切り位置情報を持つため、作成やメンテナンスの作業に多大な手間を要していた。

【0007】

そこで本発明は、形態素解析等の文章を単語に分解する処理において、複合語の分割処理を効率的に行い、処理全体における実行効率を向上させることを目的とする。

また本発明は、複合語を分割した際にも解析結果として得られるトークン列が最適なものであることを保証できるようにすることを他の目的とする。

さらに本発明は、複合語辞書の作成及びメンテナンスに要する手間を削減することをさらに他の目的とする。

【0008】

【課題を解決するための手段】

上記の目的を達成する本発明は、次のように構成された形態素解析装置として実現される。この形態素解析装置は、見出し語およびこの見出し語の属性情報を格納した辞書部と、この辞書部を参照して、処理対象の自然言語文からかかる自然言語文を構成することが可能なトークンを抽出し、トークンリストに登録するトークンリスト作成部と、このトークンリスト作成部にて作成されたトークンリストに基づいて処理対象の自然言語文を構成するのに最適なトークン列を選択するトークン列選択部とを備える。そして、トークンリスト作成部は、形態素解析に対して与えられた条件と、抽出した前記トークンに対応する見出し語の属性情報とに基づいて、このトークンのトークンリストへの登録制御を行うことを特徴とする。

この登録制御は、より詳細には、形態素解析に対して与えられた条件に合致する属性を持つトークンのみを前記トークンリストに登録することにより実現される。さらに詳細には、属性情報は、属性情報の数に応じたビット数のフラグデータで記録され、トークンリスト作成部は、このフラグデータの値に基づいて、トークンをトークンリストに登録するか否かを決定する。

【0009】

また、本発明の他の形態素解析装置は、処理対象の自然言語文をかかるとなる自然言語文の構成要素であるトークンに分解し、より小さいトークンに分割可能なトークンを除いてトークンリストに登録するトークンリスト作成手段と、このトークンリスト作成手段にて作成されたトークンリストに基づいて処理対象の自然言語文を構成するのに最適なトークン列を選択するトークン列選択手段とを備えることを特徴とする。

【0010】

上記の目的を達成する他の本発明は、形態素解析手段と、形態素解析された自然言語文に対して所定の処理を行うアプリケーション実行手段とを備えた自然言語処理装置としても実現される。この自然言語処理装置において、形態素解析手段は、見出し語およびその属性情報を格納した辞書部と、この辞書部を参照して、処理対象の自然言語文からこの自然言語文を構成することが可能なトークンを抽出し、抽出されたトークンに対応する見出し語の属性情報に基づいて、アプリケーション実行手段にて要求される条件に合致する属性を持つトークンのみをトークンリストに登録するトークンリスト作成部と、このトークンリスト作成部にて作成されたトークンリストに基づいて自然言語文を構成するのに最適なトークン列を選択するトークン列選択部とを備えることを特徴とする。アプリケーション実行手段にて実現される処理としては、例えば、テキスト検索処理、機械翻訳処理、テキスト・マイニング等が挙げられる。

【 0 0 1 1 】

さらにまた、上記の目的を達成する他の本発明は、コンピュータを用いて自然言語文に対し形態素解析を行う、次のような形態素解析方法としても実現される。この形態素解析方法は、処理対象の自然言語文を入力し、メモリに格納された辞書を参照して、この自然言語文を構成することが可能なトークン及びその属性情報を取得し、メモリの作業領域に格納するステップと、形態素解析に与えられた所定の条件およびトークンの属性情報に基づき、メモリに格納されたトークンの中から所定のトークンを選択してメモリの所定領域に構築されたトークンリストに登録するステップと、トークンリストに基づいて処理対象の自然言語文を構成することが可能なトークン列を生成し、メモリの作業領域に格納するステップと、メモリに格納されたトークン列の中で処理対象の自然言語文を構成するのに最適なトークン列を選択し出力するステップとを含むことを特徴とする。

【 0 0 1 2 】

また、本発明の他の形態素解析方法は、処理対象の自然言語文を入力し、この自然言語文の構成要素であるトークンに分解し、得られたトークン群をメモリの作業領域に格納するステップと、かかるトークン群を、より小さいトークンに分割可能なトークンを除いてメモリの所定領域に構築されたトークンリストに登録

するステップと、このトークンリストに基づいて処理対象の自然言語文を構成することが可能なトークン列を生成し、メモリの作業領域に格納するステップと、メモリに格納されたトークン列の中で処理対象の自然言語文を構成するのに最適なトークン列を選択し出力するステップとを含むことを特徴とする。

【0013】

さらに本発明は、コンピュータを制御して上述した形態素解析装置あるいは自然言語処理装置の機能を実現するプログラム、またはコンピュータに上記の形態素解析方法の各ステップに対応する処理を実行させるプログラムとしても実現される。このプログラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記録媒体に格納して配布したり、ネットワークを介して配信したりすることにより提供することができる。

【0014】

【発明の実施の形態】

以下、添付図面に示す実施の形態に基づいて、この発明を詳細に説明する。

図1は、本発明による形態素解析を実行するのに好適なコンピュータ装置のハードウェア構成の例を模式的に示した図である。

図1に示すコンピュータ装置は、演算手段であるCPU (Central Processing Unit: 中央処理装置) 101と、M/B (マザーボード) チップセット102及びCPUバスを介してCPU101に接続されたメインメモリ103と、同じくM/Bチップセット102及びAGP (Accelerated Graphics Port) を介してCPU101に接続されたビデオカード104と、PCI (Peripheral Component Interconnect) バスを介してM/Bチップセット102に接続されたハードディスク105、ネットワークインターフェイス106及びUSBポート107と、さらにこのPCIバスからブリッジ回路108及びISA (Industry Standard Architecture) バスなどの低速なバスを介してM/Bチップセット102に接続されたフロッピーディスクドライブ109及びキーボード/マウス110とを備える。

なお、図1は本実施の形態を実現するコンピュータ装置のハードウェア構成を例示するに過ぎず、本実施の形態を適用可能であれば、他の種々の構成を取るこ

とができる。例えば、ビデオカード 1 0 4 を設ける代わりに、ビデオメモリのみを搭載し、CPU 1 0 1 にてイメージデータ进行处理する構成としても良いし、ATA (AT Attachment) などのインターフェイスを介してCD-ROM (Compact Disc Read Only Memory) やDVD-ROM (Digital Versatile Disc Read Only Memory) のドライブを設けても良い。

【 0 0 1 5 】

図 2 は、本実施の形態による形態素解析エンジンの機能構成を示すブロック図である。

図 2 に示すように、本実施の形態の形態素解析エンジン 1 0 は、処理対象である文章をトークンに分解し各トークンに関するトークンリストを作成するトークンリスト作成部 1 1 と、トークンリスト作成部 1 1 が使用するマスター辞書 1 2 と、作成されたトークンリストに基づいて最適なトークン列を選択するトークン列選択部 1 3 と、トークン列選択部 1 3 が使用する文法辞書 1 4 とを備える。

【 0 0 1 6 】

上記の構成のうち、トークンリスト作成部 1 1 及びトークン列選択部 1 3 は、図 1 に示したメインメモリ 1 0 3 に展開されたプログラムにてCPU 1 0 1 を制御することにより実現される仮想的なソフトウェアブロックである。CPU 1 0 1 を制御してこれらの機能を実現させるプログラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記憶媒体に格納して配布したり、ネットワークを介して配信したりすることにより提供される。本実施の形態では、図 1 に示したネットワークインターフェイス 1 0 6 やフロッピーディスクドライブ 1 0 9、図示しないCD-ROMドライブなどを介して当該プログラムを入力し、ハードディスク 1 0 5 に格納する。そして、ハードディスク 1 0 5 に格納されたプログラムをメインメモリ 1 0 3 に読み込んで展開し、CPU 1 0 1 にて実行することにより、これらの機能を実現する。

【 0 0 1 7 】

また、マスター辞書 1 2 及び文法辞書 1 4 は、図 1 に示したメインメモリ 1 0 3 及びハードディスク 1 0 5 にて実現される。トークンリスト作成部 1 1 による処理の際にはマスター辞書 1 2 が、トークン列選択部 1 3 による処理の際には文

法辞書14がそれぞれハードディスク105からメインメモリ103に読み込まれる。そして、トークンリスト作成部11またはトークン列選択部13として機能するCPU101にて参照される。

【0018】

本実施の形態による形態素解析エンジン10は、従来の形態素解析のようにトークン列を選んだ後で複合語を分割するのではなく、トークンリストを作成する段階で複合語を考慮して処理を行う。これによって、作成されたトークンリストからトークン列の選択が行われると、複合語が分割され、かつ最適なトークン列が選択されることとなる。

なお、複合語を分割するか否かは、形態素解析の結果を使用するアプリケーションの要求に応じて選択される。例えば、文書検索やテキストマイニングでは、できるだけ多くの関連項目が検出される（ヒットする）ように、複合語を細かく分割することが好ましい場合がある。一方、機械翻訳などでは、複合語を分割してしまうと意味が変わってしまうため、複合語は分割せずに複合語のままで扱う方が好ましい場合がある。したがって、形態素解析における複合語の分割は、アプリケーションの要求に応じて選択的に（当該アプリケーションのオプション設定等に基づいて）実行される。

【0019】

上記の構成において、トークンリスト作成部11は、処理対象の文章を構成する文字列を切り出し、全ての可能性のあるトークンを得る。そして、マスター辞書12を参照してトークンリストを作成する。本実施の形態においては、トークンリスト作成部11は、複合語の分割して形態素解析を行う設定である場合、複合語に対応するトークンを除いてトークンリストに登録する。すなわち、トークンの属性に応じてトークンリストへの登録を制御する。以下、具体例を挙げて説明する。

「情報処理学会で青木和夫の」という文（の一部）に対してトークンリストを作成する場合を考える。

【0020】

図3は、この例文に対する複合語を分割しない場合のトークンリストを示す図

である。

複合語を分割しない場合、例えば「情報処理学会」という語について「情」、「情報」、「情報処理」、「情報処理学会」という文字列がそれぞれトークンとして抽出され、マスター辞書 12 から得られる品詞情報と共にトークンリストに登録される。

図 4 は、同じ例文に対する複合語を分割する場合のトークンリストを示す図である。

複合語を分割する場合、「情報処理学会」という語から切り出される文字列のうち、「情報処理」及び「情報処理学会」は複合語であるので、トークンリストに登録されない（図 3 と図 4 とを比較すると、「情報処理」、「情報処理学会」及び「青木和夫」が複合語として除去されている）。「情報処理」及び「情報処理学会」が複合語であるか否かは、後述するマスター辞書 12 に登録されている情報に基づいて判断される。

【0021】

マスター辞書 12 は、トークンとそのトークンに関する情報とが対応付けられて登録されている。

図 5 は、マスター辞書 12 におけるデータフォーマットの例を示す図である。

図 5 に示すように、マスター辞書 12 には、見出し語（トークン）ごとに、当該見出し語の品詞情報と、当該見出し語が分割可能か否かを示すフラグ（分割可能フラグ）とが登録されている（以下、見出し語自体を含むこれらの情報をトークン情報と称す）。マスター辞書の品詞情報には、正確には品詞の種類を示すものではないが、人名、地名、組織名等の属性を示す情報を含めることができる。また図示の例では、分割可能フラグの値が 0 の登録後は分割不可能であり、1 の登録後は分割可能であることを示している。トークンリスト作成部 11 は、複合語を分割する設定のときは、このフラグを参照してトークンの文字列が分割可能か否か（複合語であるか否か）を判断し、分割可能であれば当該トークンをトークンリストに登録しない。

【0022】

本実施の形態では、上述したフラグによって、トークンの文字列が分割可能な

複合語か否かを示す属性情報のみを与えているが、このフラグを拡張することにより、他の種々の属性情報をトークンに与えることができる。例えば、a、b、c、dという4つの情報を4ビットのフラグデータで表現する場合、aを1（0001）、bを2（0010）、cを4（0100）、dを8（1000）と定義すれば、複合的な属性も、abは3（0011）、bcdは14（1110）というようにビット変換して表現することができる。そして、形態素解析処理に対して与えられた条件（複合語は分割する等）に合致する属性を示すフラグの値を持つトークンのみをトークンリストに登録することができる。これにより、複合語であっても人名は分割しないというような複合的な条件にてトークンリストへの登録を制御することが可能となる。

【0023】

図6は、上記のように構成されたマスター辞書12を参照し、トークンリスト作成部11がトークンリストを作成する動作を説明するフローチャートである。

図6を参照すると、まず初期動作として、処理対象の文（以下、テキスト）が入力され、マスター辞書12の内容がハードディスク105からメインメモリ103にロードされる（ステップ601）。またこのとき、メインメモリ103に、トークンリストのための領域が確保される。なお、トークンリスト作成部11による処理の開始に先立って、分割可能な複合語の分割を行うか否かのオプション設定を行っておく。この設定は、本実施の形態の形態素解析エンジン10を利用するアプリケーションのユーザインターフェイスにおいて、設定コマンドの入力を受け付けることによって行うことができる。

【0024】

処理対象のテキストが入力されると、トークンリスト作成部11は、入力テキストの先頭の文字に着目し（ステップ602）、当該着目文字から始まる各トークンのトークン情報を順次マスター辞書12から読み出し、メインメモリ103の作業領域に格納する（ステップ603、604、605）。

例えば、上述した「情報処理学会で青木和夫の」という文を処理する場合、先頭の文字「情」が着目され、「情一名詞」、「情報一名詞」、「情報処理一名詞」、「情報処理学会一名詞」というトークン情報が読み出されることとなる。

【0025】

複合語を分解するオプション設定がオンとなっているならば、次にトークンリスト作成部11は、マスター辞書12から読み出してメインメモリ103の作業領域に保持したトークン情報の分割可能フラグを調べ、当該トークンが分割可能か否かを判断する（ステップ606、607）。当該トークンが分割不可能である場合、または複合語を分解するオプション設定がオフである場合は、当該トークン情報をメインメモリ103に用意されたトークンリストに登録する（ステップ608）。そして、ステップ604に戻り、未処理のトークンが残っているか否かを調べ、残っていれば、当該未処理のトークンに関して同様の処理（ステップ605～ステップ608）を行う。

文字「情」に着目している場合、読み出された上記4つのトークン情報「情一名詞」、「情報一名詞」、「情報処理一名詞」、「情報処理学会一名詞」が全てトークンリストに登録される。

【0026】

一方、ステップ607で、当該トークンが分割可能である場合は、当該トークンをトークンリストに登録せずにステップ604へ戻り、未処理のトークンが残っているか否かを調べる。

上記の例では、「情報処理一名詞」及び「情報処理学会一名詞」が分割可能であるので、トークンリストには「情一名詞」及び「情報一名詞」のみが登録されることとなる。

【0027】

着目文字から始まる全てのトークンに関して上記の処理（ステップ605～ステップ608）を行ったならば、次に着目文字を入力テキストの後方へ1つずらしてステップ603へ戻り（ステップ609）、同様の処理（ステップ604～ステップ608）を行う。そして、入力テキストの全ての文字を着目文字として上記の処理を完了したならば、トークンリスト作成部11による処理を終了する（ステップ603）。

上記「情報処理学会で青木和夫の」という文の場合では、文字「情」に着目して上記の処理を行い、次に文字「報」に着目して処理を行うというように、順次

処理を行っていき、最後の文字「の」に着目して処理を行った後、トークンリスト作成部 11 の処理が終了する。

【0028】

トークン列選択部 13 は、従来の形態素解析エンジンにおけるトークン列選択部と同様である。すなわち、文法辞書 14 を参照して、トークンリスト作成部 11 にて作成されたトークンリストから、最も可能性の高い（最適な）トークン列を選択する。トークン列の選択には、通常の接続コスト最小法を用いることができる。

トークン列選択部 13 の処理に利用される文法辞書 14 は、従来の形態素解析エンジンにおける文法辞書と同様である。すなわち、単語の可能なつながり方と各つながり方に対して予め設定されたコストとを定義した文法データとが格納されている。

【0029】

図 7 は、トークン列選択部 13 による処理を説明するフローチャートである。

なお、処理が開始される際、初期動作として、文法辞書 14 の内容がハードディスク 105 からメインメモリ 103 にロードされる。

図 7 に示すように、トークン列選択部 13 は、まず処理対象のテキストとトークンリスト作成部 11 にて作成されたトークンリストとを入力し（ステップ 701）、文法辞書 14 を参照して、入力テキストの先頭から末尾までを構成する可能なトークン列を生成し、メモリの作業領域に格納する（ステップ 702）。

上述した「情報処理学会で青木和夫の」という文では、例えば、

情報処理学会－で－青木和夫－の（複合語を分割しない場合）

情報－処理－学会－で－青木－和夫－の

情－報－処－理－学－会－で－青－木－和－夫－の

等のトークン列が得られる。

【0030】

次に、トークン列選択部 13 は、生成されたトークン列を解（経路）の候補として、トークン列を構成するトークンのつながり方に対して与えられたコストを文法辞書から参照し、コストの総和が最も低い最適解（最短経路）を探索する（

ステップ 7 0 3)。この探索には、例えば周知のダイクストラ・アルゴリズムを用いることができる。

最後に、トークン列選択部 1 3 は、最適解であるトークン列を、入力テキストに対する最適なトークン列として出力する（ステップ 7 0 4）。

【 0 0 3 1 】

上述したトークン列選択部 1 3 による処理は、従来の形態素解析エンジンにおける処理と同様である。しかしながら、分割可能な複合語を分割する設定で形態素解析を行う場合、上述したようにトークンリスト作成部 1 1 によるトークンリストの作成処理の段階で不要な複合語がトークンリストから除去されているため、その分だけ処理対象となるトークンの組み合わせの数（パスの数）も、従来の形態素解析の場合と比べて少なくなる。したがって、トークン列選択部 1 3 による処理が高速化されることとなる。

【 0 0 3 2 】

また、従来の形態素解析エンジンでは、分割可能な複合語を分割する設定で形態素解析を行う場合、トークン列選択部 1 3 によって選択されたトークン列に対して複合語辞書とのマッチングを行い、当該トークン列に含まれる分割可能な複合語を分割していた。そのため、複合語辞書を持つ分だけメモリやハードディスクといった記憶装置（資源）の使用量が増すと共に、形態素解析処理の実行時に複合語を分割するための時間が余計にかかっていた。

これに対し、本実施の形態の形態素解析エンジン 1 0 では、分割可能な複合語を分割する設定で形態素解析を行う場合、トークンリスト作成部 1 1 にて作成されたトークンリストには複合語のトークンを登録しないことによって複合語の分割に対応するため、マスター辞書の他に複合語辞書を持つ必要が無く、記憶装置（資源）使用量を削減できる。そして、形態素解析処理の実行時にも、トークンリストの作成およびトークン列の選択の他に複合語の分割処理を行う必要がないため、処理に要する時間を短縮できる。

【 0 0 3 3 】

さらに、従来の形態素解析エンジンは、上記のように最適なトークン列を選択した後に分割可能な複合語を分割していたため、当該トークン列が最適なのは複

合語を複合語のまま扱う場合であり、当該複合語が分割された状態のトークン列が最適であることは保証されない。

これに対し、本実施の形態素解析エンジン 10 では、分割可能な複合語を分割する設定で形態素解析を行う場合、トークンリスト作成部 11 にて作成されたトークンリストには複合語のトークンが含まれておらず、複合語のトークンを含むトークン列がトークン列選択部 13 による処理対象となることはない。したがって、トークン列選択部 13 により選択されたトークン列は必ず複合語のトークンを含まないトークン列であり、かつ最適であることが保証される。

【0034】

次に、本実施の形態の形態素解析エンジン 10 が利用されるアプリケーションについて説明する。

上述した形態素解析エンジン 10 は、コンピュータ装置に搭載されて自然言語文に対する形態素解析装置として実現される他、テキスト検索システムや機械翻訳システム、テキスト・マイニング・システム等の自然言語に対する処理を行う種々のアプリケーションにて利用することができる。

図 8 は、形態素解析エンジン 10 を搭載したテキスト検索システムの構成例を示す図である。

図 8 を参照すると、このテキスト検索システムは、検索対象のテキスト群を格納したテキストデータベース 801 と、テキストデータベース 801 に格納されている各テキストからキーワードのインデックスファイルを作成するインデックスファイル作成部 802 と、インデックスファイルを用いて検索対象のテキスト群に対し検索タームである文の検索を行うテキスト検索部 803 と、インデックスファイル作成部 802 及びテキスト検索部 803 の処理の前処理として形態素解析を行う形態素解析部 804 と、検索タームである文を入力するテキスト入力部 805 と、検索結果を出力する検索結果出力部 806 とを備える。

【0035】

このテキスト検索システムは、例えば 1 台またはネットワークで接続された複数台のコンピュータ装置にて実現される。図 8 に示した構成において、テキストデータベース 801 はハードディスク等の記憶手段にて実現され、テキスト検索

というアプリケーションを実行する手段であるインデックスファイル作成部802及びテキスト検索部803はプログラム制御されたCPUにて実現される。また、形態素解析部804として本実施の形態の形態素解析エンジン10を用いることができる。テキスト入力部805はキーボードやマウス、その他の入力デバイスで実現され、検索結果出力部806はディスプレイ装置等で実現される。また、ネットワークインターフェイスを介して外部機器との間で検索タームである文の入力や検索結果の出力を行っても良い。

【0036】

このテキスト検索システムでは、インデックスファイルの作成時とテキスト検索の実行時とに形態素解析が行われる。

インデックスファイルの作成処理において、まずテキストデータベース801からテキストが順次読み出され、形態素解析部804により形態素解析が行われる。このとき、テキスト検索（アプリケーション）の必要から複合語を分割したい場合は、図6、図7に示したように、複合語を含まないトークン列の中から最適なものが選択される。得られたトークン列の中から、インデックスファイル作成部802により、名詞や形容詞等の自立語のトークン（単語）のみがキーワードとして抜き出される。そして、検索対象のテキストごとにキーワードが登録されたインデックスファイルが作成される。

【0037】

次にテキスト検索の処理において、まずテキスト入力部805により検索タームである文が入力され、形態素解析部804により当該入力文の形態素解析が行われる。このとき、テキスト検索（アプリケーション）の必要から複合語を分割したい場合は、図6、図7に示したように、複合語を含まないトークン列の中から最適なものが選択される。得られたトークン列の中から、テキスト検索部803により、名詞や形容詞等の自立語のトークン（単語）が抜き出される。そして、インデックスファイルを用いて、抜き出されたトークンを含むテキストの検索が行われる。この検索によりヒットしたテキストがテキストデータベース801から読み出され、検索結果出力部806にて出力（表示）される。

【0038】

図 9 は、形態素解析エンジン 1 0 を搭載した機械翻訳システムの構成例を示す図である。

図 9 を参照すると、この機械翻訳システムは、翻訳対象である原文テキストを入力する原文入力部 9 0 1 と、入力された原文テキストを形態素解析する形態素解析部 9 0 2 と、形態素解析の行われた原文テキストを構文解析する構文解析部 9 0 3 と、構文解析の結果に基づいて原文テキストの文構造を翻訳文の言語における文構造に構文変換する構文変換部 9 0 4 と、構文変換の結果得られた文構造に基づいて翻訳文テキストを生成する翻訳文生成部 9 0 5 と、生成された翻訳文テキストを出力する翻訳文出力部 9 0 6 とを備える。また、特に図示しないが、原文及び翻訳文の各言語における単語辞書や文法辞書を備え、各部の処理において用いられる。

【 0 0 3 9 】

この機械翻訳システムは、例えば 1 台またはネットワークで接続された複数台のコンピュータ装置にて実現される。図 9 に示した構成において、形態素解析部 9 0 2 として本実施の形態の形態素解析エンジン 1 0 を用いることができる。機械翻訳というアプリケーションを実行する手段である構文解析部 9 0 3、構文変換部 9 0 4 及び翻訳文生成部 9 0 5 は、プログラム制御された CPU にて実現される。また、原文入力部 9 0 1 はキーボードやマウス、その他の入力デバイスで実現され、翻訳文出力部 9 0 6 はディスプレイ装置等で実現される。また、ネットワークインターフェイスを介して外部機器との間で原文テキストの入力や翻訳文テキストの出力を行っても良い。

【 0 0 4 0 】

機械翻訳では、複合語を分割するか否か等、単語の属性に応じて訳出の仕方を調整することによって翻訳の精度が大きく変わるため、詳細な設定が可能であることが好ましい。本実施の形態の形態素解析エンジンを用いた形態素解析部 9 0 2 によれば、そのような設定に応じて、トークンリストの作成段階で不要な（設定に合致しない属性を持つ）トークンを排除し、残りのトークンから最適なトークン列を得ることができる。

【 0 0 4 1 】

図 1 0 は、形態素解析エンジン 1 0 を搭載したテキスト・マイニング・システムの構成例を示す図である。

図 1 0 を参照すると、このテキスト・マイニング・システムは、分類対象のテキスト群を格納したテキストデータベース 1 0 0 1 と、テキストの分類基準となる分類表を格納した分類表格納部 1 0 0 2 と、分類表を参照してテキストデータベース 1 0 0 1 に格納されている各テキストの分類を行うマッチング処理部 1 0 0 3 及び分類実行部 1 0 0 4 と、マッチング処理部 1 0 0 3 の処理の前処理として形態素解析を行う形態素解析部 1 0 0 5 と、分類されたテキストを格納する分類テキストデータベース 1 0 0 6 とを備える。

【 0 0 4 2 】

このテキスト・マイニング・システムは、例えば 1 台またはネットワークで接続された複数台のコンピュータ装置にて実現される。図 1 0 に示した構成において、テキストデータベース 1 0 0 1、分類表格納部 1 0 0 2 及び分類テキストデータベース 1 0 0 6 はハードディスク等の記憶手段にて実現され、テキスト・マイニングというアプリケーションを実行する手段であるマッチング処理部 1 0 0 3 及び分類実行部 1 0 0 4 はプログラム制御された CPU にて実現される。また、形態素解析部 1 0 0 5 として本実施の形態の形態素解析エンジン 1 0 を用いることができる。

【 0 0 4 3 】

このテキスト・マイニング・システムでは、まずテキストデータベース 1 0 0 1 からテキストが順次読み出され、形態素解析部 1 0 0 5 により形態素解析が行われる。このとき、テキスト・マイニング（アプリケーション）の必要から複合語を分割したい場合は、図 6、図 7 に示したように、複合語を含まないトークン列の中から最適なものが選択される。得られたトークン列の中から、マッチング処理部 1 0 0 3 により、名詞や形容詞等の自立語のトークン（単語）のみがキーワードとして抜き出される。そして、分類表格納部 1 0 0 2 から単語と当該単語を含むテキストのカテゴリとを対応付けて登録した分類表が読み出され、予め定められた所定のルールに従って、トークン列から抽出された単語と分類表の単語とのマッチングが行われる。

次に、マッチング処理部 1 0 0 3 によるマッチングの結果（単語の割合等）に基づいて、分類実行部 1 0 0 4 により、着目中のテキストのカテゴリが決定される。そして、決定されたカテゴリに応じて分類テキストデータベース 1 0 0 6 に格納される。

【 0 0 4 4 】

なお、上記実施の形態では、日本語や中国語、韓国語などのように単語を区切らない表記方法を採用言語にて記述された自然言語文を解析する際に行われる形態素解析について説明したが、その他の言語においても、接頭辞や接尾辞、その他の複合語を適切に分割することで自然言語文の処理の性能向上を期待することができる場合に、本発明を適用することが可能である。

【 0 0 4 5 】

【発明の効果】

以上説明したように、本発明によれば、形態素解析等の文章を単語に分解する処理において、複合語の分割処理を効率的に行い、処理全体における実行効率の向上を図ることができる。

また本発明によれば、複合語を分割した際にも解析結果として得られるトークン列が最適なものであることを保証することができる。

さらに本発明によれば、複合語辞書の作成及びメンテナンスに要する手間を削減することができる。

【図面の簡単な説明】

【図 1】 本実施の形態による形態素解析を実行するのに好適なコンピュータ装置のハードウェア構成の例を模式的に示した図である。

【図 2】 本実施の形態による形態素解析エンジンの機能構成を示す図である。

【図 3】 複合語を分割しない場合のトークンリストの例を示す図である。

【図 4】 図 3 と同じ例文に対する、本実施の形態による複合語を分割する場合のトークンリストの例を示す図である。

【図 5】 本実施の形態におけるマスター辞書におけるデータフォーマットの例を示す図である。

【図 6】 本実施の形態におけるトークンリストを作成する動作を説明するフローチャートである。

【図 7】 本実施の形態におけるトークン列選択部による処理を説明するフローチャートである。

【図 8】 本実施の形態の形態素解析エンジンを搭載したテキスト検索システムの構成例を示す図である。

【図 9】 本実施の形態の形態素解析エンジンを搭載した機械翻訳システムの構成例を示す図である。

【図 10】 本実施の形態の形態素解析エンジンを搭載したテキスト・マイニング・システムの構成例を示す図である。

【図 11】 コンピュータにて実現される従来の形態素解析手段の機能ブロックを示す図である。

【図 12】 従来の形態素解析のアルゴリズムを概略的に説明するフローチャートである。

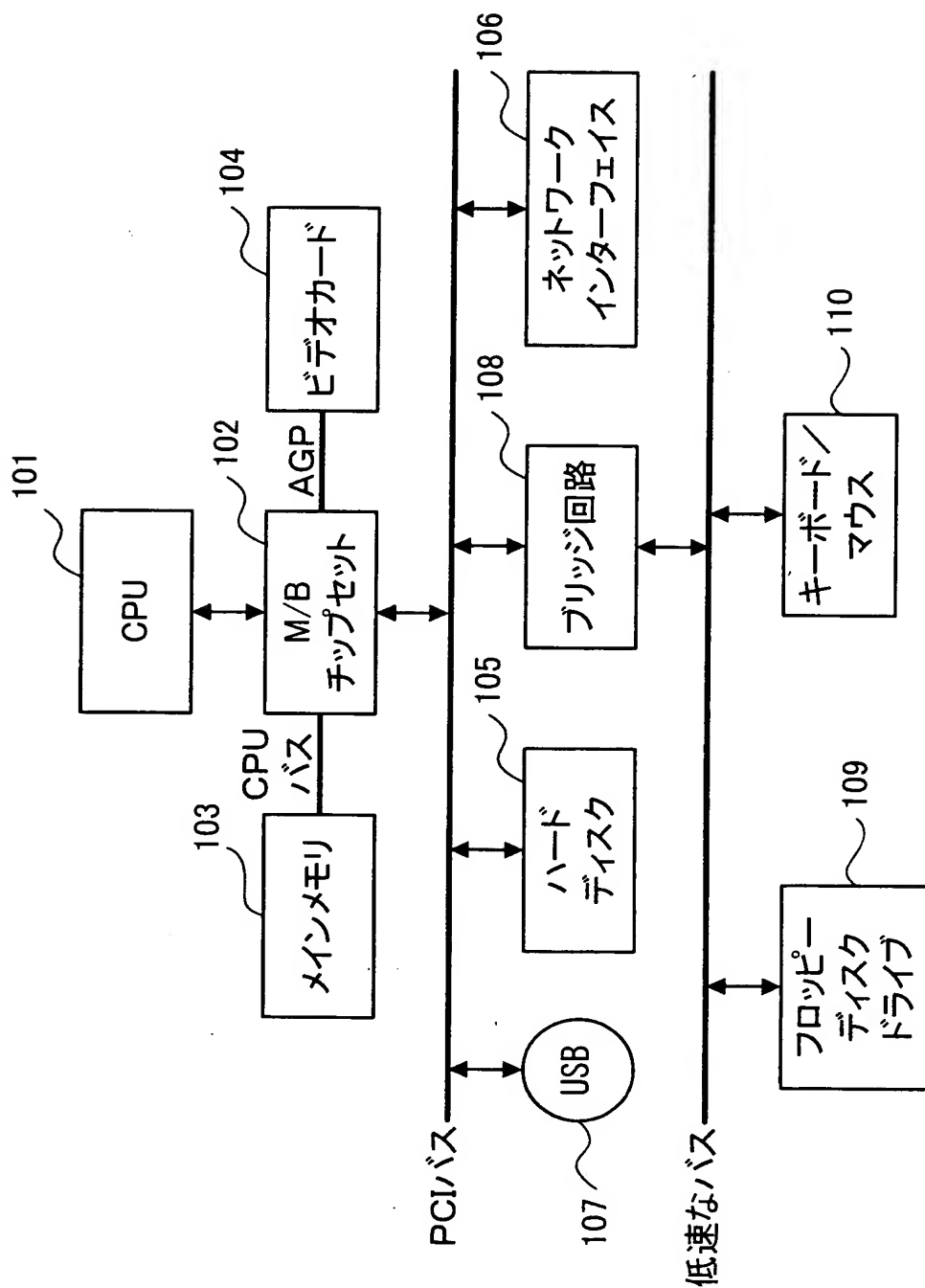
【符号の説明】

1 0 …形態素解析エンジン、1 1 …トークンリスト作成部、1 2 …マスター辞書、1 3 …トークン列選択部、1 4 …文法辞書、1 0 1 …CPU、1 0 2 …M/Bチップセット、1 0 3 …メインメモリ、1 0 5 …ハードディスク、1 0 6 …ネットワークインターフェイス、8 0 1、1 0 0 1 …テキストデータベース、8 0 2 …インデックスファイル作成部、8 0 3 …テキスト検索部、8 0 4、9 0 2、1 0 0 5 …形態素解析部、8 0 5 …テキスト入力部、8 0 6 …検索結果出力部、9 0 1 …原文入力部、9 0 3 …構文解析部、9 0 4 …構文変換部、9 0 5 …翻訳文生成部、9 0 6 …翻訳文出力部、1 0 0 2 …分類表格納部、1 0 0 3 …マッチング処理部、1 0 0 4 …分類実行部、1 0 0 6 …分類テキストデータベース

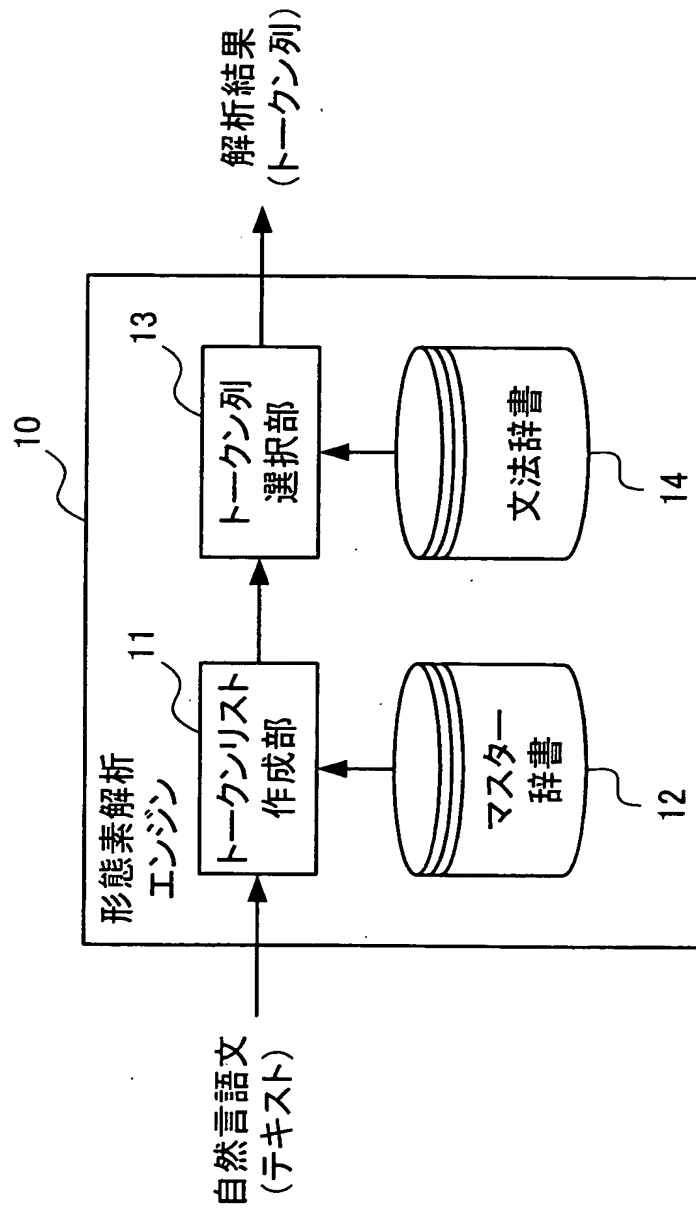
【書類名】

図面

【図 1】



【図2】



【図 3】

情	名詞	青	五段動詞
情報	名詞	青	形容詞
情報処理	名詞	青	名詞
情報処理学会	名詞	青	固有名詞一姓
		青木	固有名詞一姓
		青木	固有名詞一名
		青木	固有名詞一地名
		青木和夫	固有名詞一姓名

報	五段動詞	木	
報	一段動詞		:
報	サ変動詞		
報	名詞		
報	接尾辞		

処	五段動詞		
処	サ変動詞		
処理	サ変動詞		
処理	名詞		

理	名詞		
理	固有名詞一名前		
理学	名詞		

学	五段動詞		
学	名詞		
学	固有名詞一名前		
学会	名詞		

会	五段動詞		
会	サ変動詞		
会	接尾辞		

で	Unknown		
---	---------	--	--

		の	Unknown
--	--	---	---------

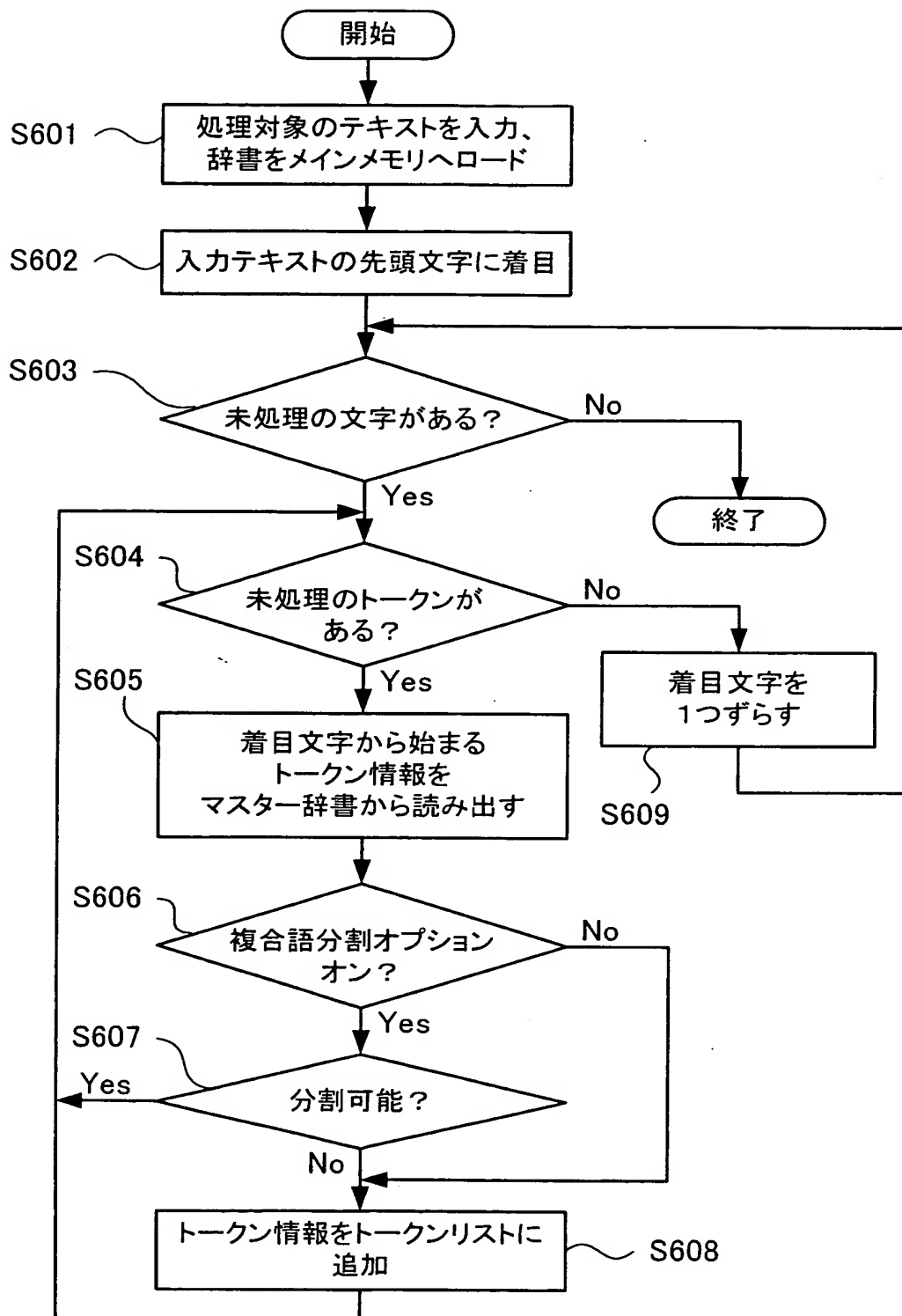
【図 4】

情 情報	名詞 名詞	青 青 青 青 青木 青木 青木	五段動詞 形容詞 名詞 固有名詞一姓 固有名詞一姓 固有名詞一名 固有名詞一地名
報 報 報 報 報	五段動詞 一段動詞 サ変動詞 名詞 接尾辞	木	:
処 処 処理 処理	五段動詞 サ変動詞 サ変動詞 名詞	和	:
理 理 理学	名詞 固有名詞一名前 名詞	夫	:
学 学 学 学会	五段動詞 名詞 固有名詞一名前 名詞	の	Unknown
会 会 会	五段動詞 サ変動詞 接尾辞		
で	Unknown		

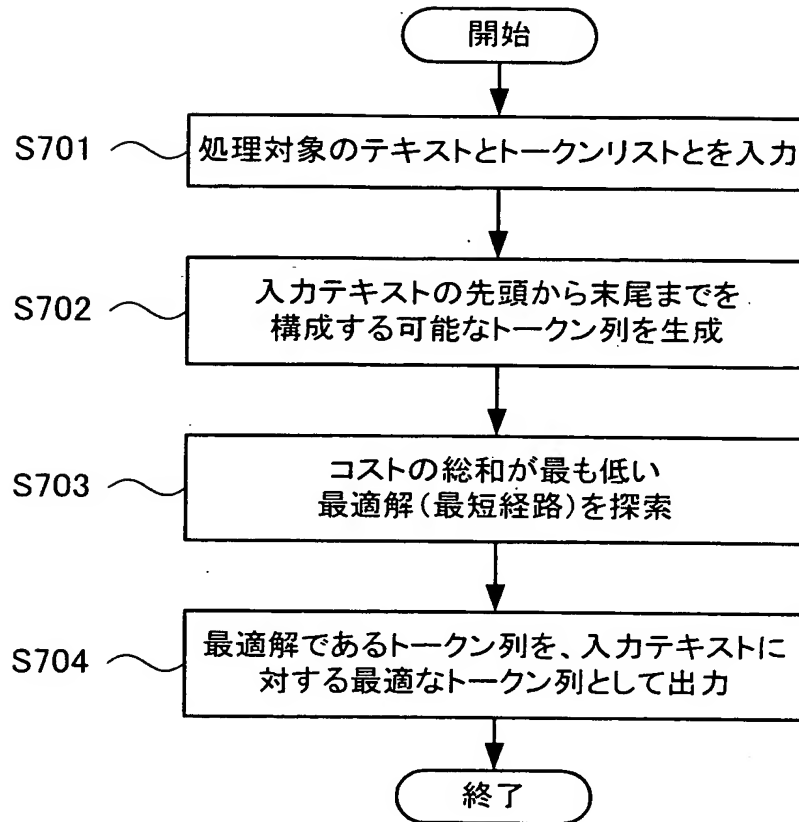
【図 5】

登録後	品詞情報	分割可能フラグ
情報	名詞	0
情報処理	名詞	1
情報処理学会	名詞	1
⋮	⋮	⋮

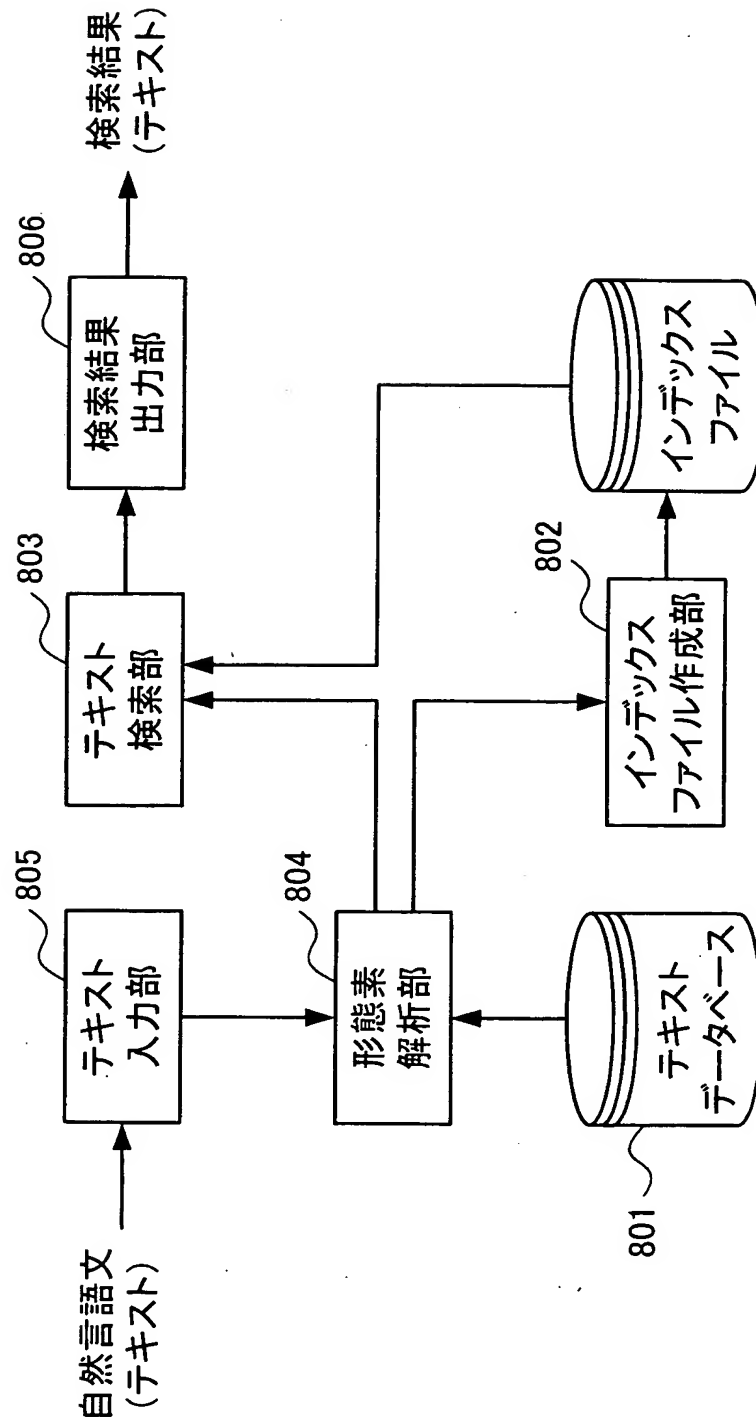
【図 6】



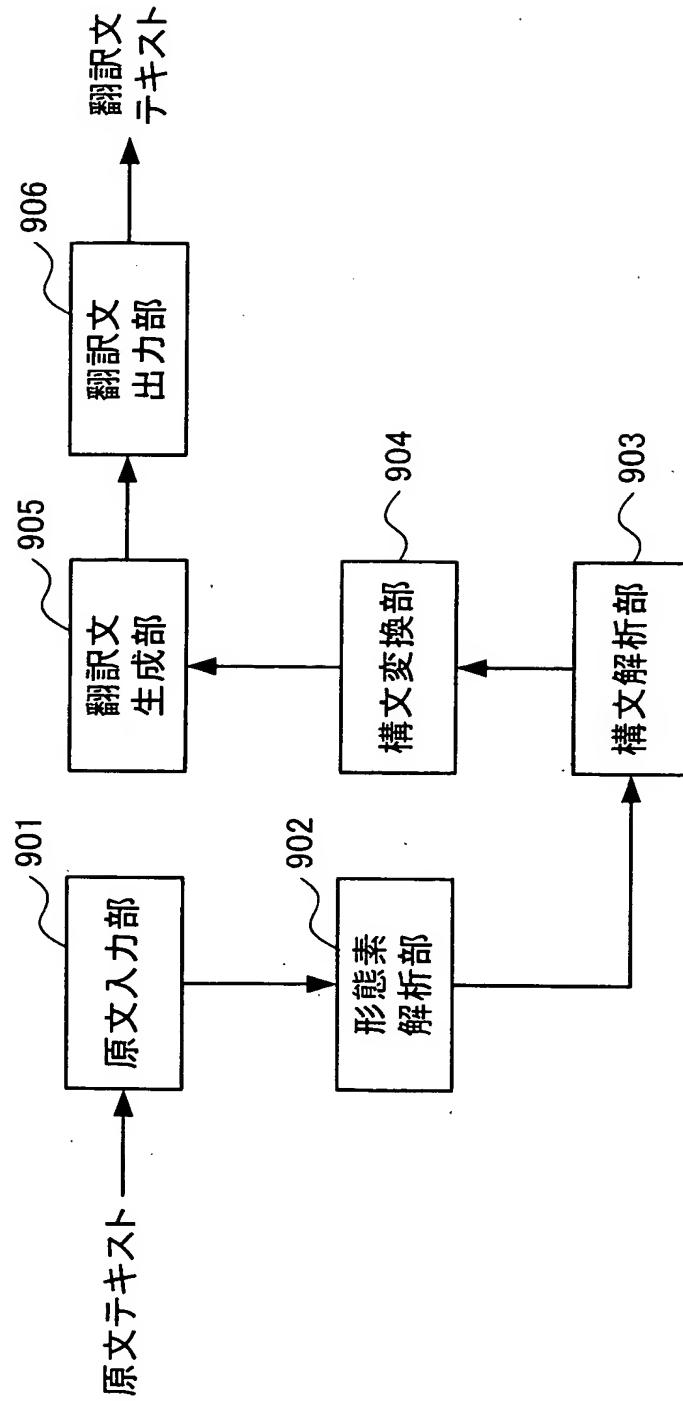
【図 7】



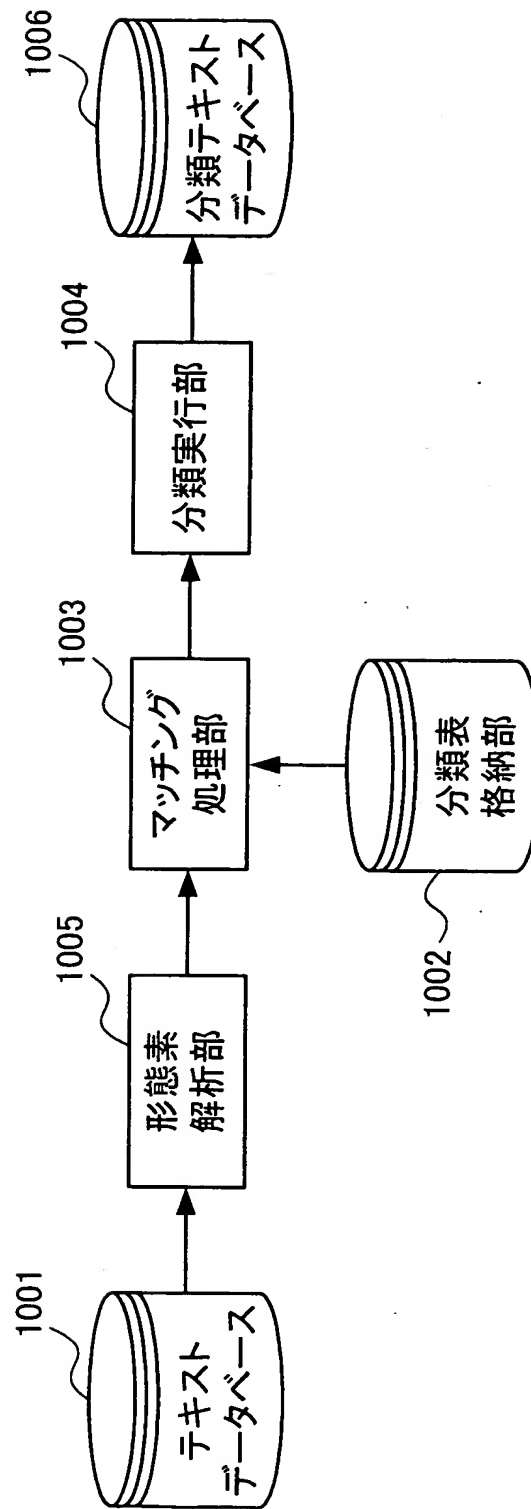
【図 8】



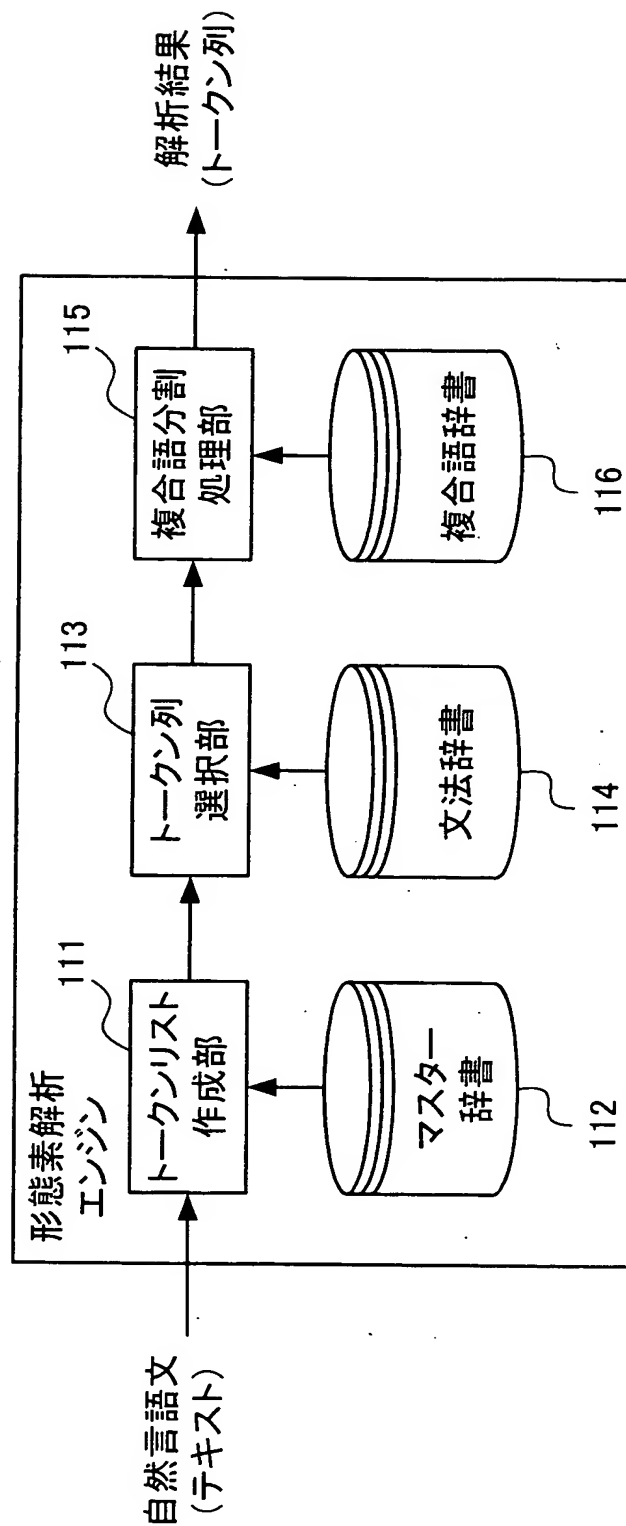
【図 9】



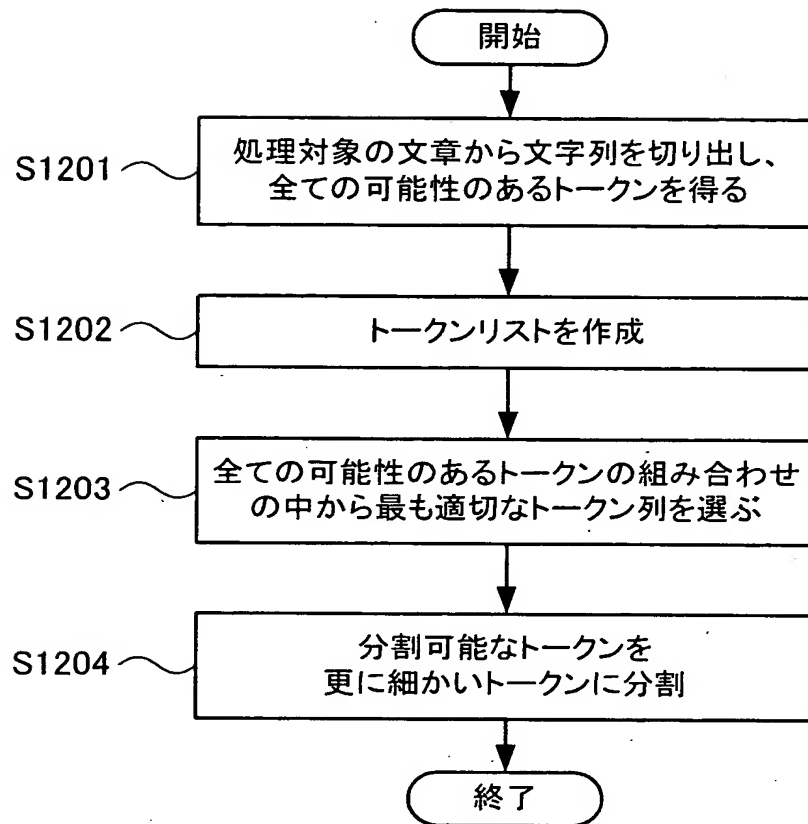
【図 10】



【図 11】



【図 1 2】



【書類名】 要約書

【要約】

【課題】 形態素解析等の文章を単語に分解する処理において、複合語の分割処理を効率的に行い、かつ複合語を分割した際にも解析結果として得られるトークン列が最適なものであることを保証できるようにする。

【解決手段】 処理対象の自然言語文をかかる自然言語文の構成要素であるトークンに分解してトークンリストに登録するトークンリスト作成部 1 1 と、このトークンリスト作成部 1 1 にて作成されたトークンリストに基づいて処理対象の自然言語文を構成するのに最適なトークン列を選択するトークン列選択部 1 3 とを備える。そして、トークンリスト作成部 1 1 は、形態素解析に対して与えられた条件に応じて、処理対象の自然言語文を分解して得られたトークンのうち、より小さいトークンに分割可能なトークンを除いてトークンリストに登録する。

【選択図】 図 2

認定・付加情報

特許出願の番号	特願 2003-033220
受付番号	50300215595
書類名	特許願
担当官	土井 恵子 4264
作成日	平成 15 年 3 月 25 日

<認定情報・付加情報>

【特許出願人】

【識別番号】	390009531
【住所又は居所】	アメリカ合衆国 10504、ニューヨーク州 アーモンク ニュー オーチャード ロード
【氏名又は名称】	インターナショナル・ビジネス・マシーンズ・コーポレーション

【代理人】

【識別番号】	100086243
【住所又は居所】	神奈川県大和市下鶴間 1623 番地 14 日本アイ・ビー・エム株式会社 大和事業所内
【氏名又は名称】	坂口 博

【代理人】

【識別番号】	100091568
【住所又は居所】	神奈川県大和市下鶴間 1623 番地 14 日本アイ・ビー・エム株式会社 大和事業所内
【氏名又は名称】	市位 嘉宏

【代理人】

【識別番号】	100108501
【住所又は居所】	神奈川県大和市下鶴間 1623 番 14 日本アイ・ビー・エム株式会社 知的所有権
【氏名又は名称】	上野 剛史

【復代理人】

【識別番号】	100104880
【住所又は居所】	東京都港区赤坂 5-4-11 山口建設第 2 ビル 6 F セリオ国際特許事務所
【氏名又は名称】	古部 次郎

【選任した復代理人】

【識別番号】	100118201
--------	-----------

次頁有

認定・付加情報（続き）

【住所又は居所】 東京都港区赤坂5-4-11 山口建設第二ビル
6F セリオ国際特許事務所
【氏名又は名称】 千田 武

出 願 人 履 歴 情 報

識別番号 [390009531]

1. 変更年月日 2002年 6月 3日

[変更理由] 住所変更

住 所 アメリカ合衆国10504、ニューヨーク州 アーモンク ニ
ュー オーチャード ロード

氏 名 インターナショナル・ビジネス・マシーンズ・コーポレーショ
ン